

### **Amendments to the Specification:**

In the following amendments, insertions are double-underlined and deletions are marked with strikethrough.

On page 1, line 1, of the specification, please replace the paragraph that was added in the Preliminary Amendment filed January 28, 2002, with the following amended paragraph:

#### **CROSS REFERENCE TO RELATED APPLICATION**

This application is a continuation of United States Application no. 09/222,596, filed December 28, 1998, now U.S. Patent No. 6,351,712 B1, which is incorporated by reference herein in its entirety.

On page 3, line 14, of the specification, please replace the paragraph beginning “The use of two fluorophores has been described” with the following amended paragraph:

The use of two fluorophores has been described by Shalon *et al.* ~~Shalon et al.~~, 1996, “A microarray system for analyzing complex DNA samples using two-color fluorescent probe hybridization,” *Genome Research* 6:629-645. The problem with the approach put forth by Shalon *et al.* is that each species of mRNA molecule has a bias in its measured color ratio due to interaction of the fluorescent labeling molecule with either the reverse transcription of the mRNA or with the hybridization efficiency or both. Without any error correction scheme to account for this bias, the data from a single microarray experiment, or even a plurality of nominal repeats of a microarray experiment in which the various results are averages, will produce an unacceptable error rate. As used herein, the term nominal repeat or nominally repeated experiment refers to experiments that are run under essentially the same or similar experimental conditions such that it would be useful to combine the results of the repeated experiments.

On page 9, line 22, of the specification, please replace the paragraph beginning “Biological System: As used herein” with the following amended paragraph:

Biological System: As used herein, the term “biological system” is broadly defined to include any cell, tissue, organ or multicellular organism. For example, a biological system can be a cell line, a cell culture, a tissue sample obtained from a subject, a *Homo sapien*, a mammal,

a yeast substantially isogenic to *Saccharomyces cerevisiae*, or any other art recognized biological system. The state of a biological system can be measured by the content, activities or structures of its cellular constituents. The state of a biological system, as used herein, is determined by the state of a collection of cellular constituents, which are sufficient to characterize the cell or organism for an intended purpose including characterizing the effects of a drug or other perturbation. The term “cellular constituent” encompasses any kind of measurable biological variable. The measurements and/or observations made on the state of these constituents can be of their abundances (*i.e.*, amounts or concentrations in a biological system), their activities, their states of modification (*e.g.*, phosphorylation), or other art recognized measurements relevant to the physiological state of a biological system. In various embodiments, this invention includes making such measurements and/or observations on different collections of cellular constituents. These different collections of cellular constituents are also called aspects of the biological state of a biological system.

On page 11, line 5, of the specification, please replace the paragraph beginning “Quantitative Measurement” with the following amended paragraph:

Quantitative Measurement of Cellular Constituents: Microarrays Determining the relative abundance of diverse individual sequences in complex DNA samples is often accomplished using microarrays. *See e.g.* Shalon *et al.*, 1996, “A Microarray System for Analyzing Complex Samples Using Two-color Fluorescent Probe ~~Hybridization~~, Hybridization,” *Genome Research* 6:639-645). Frequently, transcript arrays are produced by hybridizing detectably labeled polynucleotides representing the mRNA transcripts present in a cell (*e.g.*, fluorescently labeled cDNA synthesized from total cell mRNA) to a microarray. A microarray is a surface with an ordered array of binding (*e.g.*, hybridization) sites for products of many of the genes in the genome of a cell or organism, preferably most or almost all of the genes. Microarrays are highly reproducible and therefore multiple copies of a given array can be produced and the nominal copies can be compared with each other. Preferably microarrays are small, usually smaller than 5 cm<sup>2</sup>, and made from materials that are stable under binding (*e.g.*, nucleic acid hybridization) conditions. A given binding site or unique set of binding sites in the microarray will specifically bind the product of a single gene in the cell.

On page, 14, line 16, of the specification, please replace the paragraph beginning “As detailed in the background section” with the following amended paragraph:

As detailed in the background section the two-color fluorescent hybridization process put forth by Shalon ~~et al.~~ *et al.*, *supra*, introduces bias into the profile analysis because each species of mRNA that is labeled with fluorophore has a bias in its measured color ratio due to interaction of the fluorescent labeling molecule (fluorophore) with either the reverse transcription of the mRNA or with the hybridization efficiency or both. This bias can be illustrated using the following equations. If we represent the actual molecular abundance of a particular species of mRNA  $k$ , representing cellular constituent or gene  $k$  in the biological system of interest, as  $a(k)$ , the color ratio for probe  $k$ , ignoring any source of fluorophore bias may be represented as:

$$r_{X/Y} = a_1(k) / a_2(k) \quad (1)$$

where

the subscripts 1 and 2 refer to two independently extracted mRNA cultures in which abundances are being compared;

$a_1(k)$  is the abundance of species  $k$  in mRNA culture 1;

$a_2(k)$  is the abundance of species  $k$  in mRNA culture 2;

subscripts X and Y represent the two different fluorescent labels used; and

$r_{X/Y}$  is the color ratio that ideally reflects abundance ratio  $a_1/a_2$ .

Equation (1) ideally represents the measurement plotted on the vertical axis of Figures 2 thru 6. However the use of a fluorophore labeled deoxynucleotide triphosphates affects the efficiency by which mRNA is reverse transcribed into cDNA and affects the efficiency to which the ~~fluorophore-labeled~~ fluorophore-labeled cDNA hybridizes to the microarray. The precise amount a specific fluorophore affects the transcription or hybridization efficiency is highly dependent upon the precise molecular structure of the fluorophore used. Thus, a direct comparison of  $a_1(k)$  to  $a_2(k)$ , when  $a_1(k)$  and  $a_2(k)$  are determined using different fluorophores, does not account for these fluorophore-specific affects on transcription and hybridization efficiency. The efficiency of a scanner at determining the abundances  $a_1(k)$  and  $a_2(k)$  on a microarray is also fluorophore specific. If we represent the combined efficiencies of particular fluorophore in extraction, labeling, reverse transcription, hybridization, and optical scanning as  $E$ , a more realistic representation of the color ratio presented in Equation 1 is:

$$r_{X/Y} = a_1(\mathbf{k})E_X(\mathbf{k}) / a_2(\mathbf{k})E_Y(\mathbf{k}) \quad (2)$$

where

$r_{X/Y}$  is color ratio;

the subscripts 1 and 2 are as defined for equation 1;

$a_1(\mathbf{k})$  and  $a_2(\mathbf{k})$  are as defined for equation 1;

subscripts X and Y are two fluorescent labels;

$E_X(\mathbf{k})$  is the efficiency of fluorescent ~~fluorescent~~ label X; and

$E_Y(\mathbf{k})$  is the efficiency of fluorescent ~~fluorescent~~ label Y.

In equation 2, Culture 1 has been analyzed using fluorophore X whereas Culture 2 has been analyzed using fluorophore Y. Now the color ratio  $r$  is related to the desired abundance ratio  $a_1/a_2$  but includes a factor due to the fluorophore specific efficiency biases. If a second hybridization experiment is performed, wherein Culture 1 is now analyzed with fluorophore Y and Culture 2 is analyzed using fluorophore X, the color ratio in the second hybridization experiment may be represented as:

$$r_{X/Y}^{(\text{rev})} = a_2(\mathbf{k})E_X(\mathbf{k}) / a_1(\mathbf{k})E_Y(\mathbf{k}) \quad (3)$$

where

$r_{X/Y}^{(\text{rev})}$  is color ratio in the reverse experiment; and

$a_2(\mathbf{k})$ ,  $a_1(\mathbf{k})$ ,  $E_X(\mathbf{k})$ , and  $E_Y(\mathbf{k})$  are as described for equation (2).

Performing hybridization experiments in pairs, with the label assignment reversed in one member of the pair, allows for creation of a combined average measurement in which the fluorophore specific bias is sharply reduced. For example a pair of ~~two-fluorophore~~ two-fluorophore hybridization experiments may be performed. The first two-fluorophore experiment would be performed in accordance with equation (2) and the second two-fluorophore hybridization experiments would be performed according to equation (3). If the log of the ratio of the two experiments is taken, the combined experiment can be expressed as:

$$\begin{aligned} (1/2)(\log(r_{X/Y}) - (\log(r_{X/Y}^{\text{rev}}))) &= \log(a_1(\mathbf{k})/a_2(\mathbf{k})) + (\log(E_X(\mathbf{k})/E_Y(\mathbf{k})) - \log(E_X(\mathbf{k})/E_Y(\mathbf{k}))) \\ &= \log(a_1(\mathbf{k})/a_2(\mathbf{k})) \end{aligned} \quad (4)$$

which is the desired log abundance ratio. Cancellation of the bias terms  $\log(E_X(k)/E_Y(k))$  and  $\log(E_X(k)/E_Y(k))$  relies on constancy of the biases between the first and second hybridization experiments in each fluorophore-reversed pair. Equation (4) can be written equivalently using ratios as found in equations (1)-(3) instead of differences of log ratios. However, changes in constituent levels are most appropriately expressed as the logarithm of the ratio of abundance in the pair of conditions forming the differential measurement. This is because fold changes are more meaningful than changes in absolute level, biologically.

On page 16, line 14, of the specification, please replace the paragraph beginning “This method of bias removal is particularly useful” with the following amended paragraph:

This method of bias removal is particularly useful in two-color hybridization experiments. Figure 4 illustrates the bias removal method of the present invention. Figure 4a is a color ratio vs. intensity plot for a two-color hybridization experiment in which the two cultures used are nominally the same background strain of the yeast ~~*S. Cerevisiae*~~ *S. cerevisiae*. Because the two cultures are nominally the same, it is expected that individual spots on the microarray would ~~fluoresce~~ fluoresce with the same amount of intensity for both of the fluorophores used. Experimental methods are described in the experimental section *infra*. However, as is readily apparent from Figure 4a, some of the spots on the microarray exhibit fluorophore-specific intensity. For example, spots on the microarray, corresponding to various genes in the yeast ~~*S. Cerevisiae*~~ *S. cerevisiae*, in which the intensity of the ‘red’ fluorophore is a factor of 2 ~~two~~ or more greater than the corresponding ‘green’ intensity are flagged because of their strong ~~fluorophore-specific~~ fluorophore-specific bias. Figure 4b shows the result of the fluorophore-reversed version of the experiment plotted in Figure 4a. The flagged genes in Figure 4b now have opposite bias. Figure 4c shows the result of combining the data of Figures 4a and 4b according to the methods of the present invention described above. The biases of the flagged genes have been greatly reduced.

On page 17, line 13, of the specification, please replace the paragraph beginning “If a gene of interest is present in the top 5%” on with the following amended paragraph:

If a gene of interest is present in the top 5% of up regulations in a first and second nominal repeat of a microarray experiment, the chance that it appeared that up regulated by

chance in both arrays is only  $0.05 * 0.05 = .0025$  or .25%, assuming systematic biases have been removed. Thus repeating the measurement allows a much higher level of confidence in declaring that the gene of interest is up regulated. In general, if expression ratios in any number of repeated experiments are expressed as percentile rankings, the chance  $P(H_0)$  that any (pre-specified) gene of interest is not actually up regulated is

$$P(H_0^+) = \prod_i P_i \quad (5)$$

where  $P_i$  is the percentile rank in the  $i^{\text{th}}$  experiment, expressed as a fraction (fifth percentile = 0.05). The probability that the gene is not *down*-regulated is given by

$$P(H_0^-) = \prod_i (1 - P_i) \quad (6)$$

These rank-based methods provide a powerful way of reducing false alarms with repeated measurements. For example, setting a threshold at the upper 5% of expression ratios in a hybridization to probes covering the yeast genome, which has approximately 6000 genes, would yield  ~~$6000 * 0.05 = 300$~~   $\sim 6000 * 0.05 = 300$  false detections in a single experiment, but less than one false detection on average if the same 5% threshold were applied across four experiment repeats ( $6000 * (0.05)^4$ ). This rank combining has the advantage that it does not require any modeling of the detailed error behavior in the underlying hybridization experiments, other than the assumption of no systematic biases. The rank based method is an example of a non-parametric statistical test for the significance of observed up- or down- regulations.

On page 17, line 36, of the specification, please replace the paragraph beginning “Percentile rankings such as equations (5) and (6)” with the following amended paragraph:

Percentile rankings such as equations (5) and (6) are based upon the assumption that the underlying error behavior is similar for all genes. This is not necessarily the case. For example, in Figure 5, which plots the expression ratio of two nominative repeats of the same experiment, the weakly expressing genes, as expressed by  $\log_{10}(\text{intensity})$ , have a  $\log_{10}(\text{expression ratio})$  that ~~deviate~~ deviates from the ideal value of zero. Further, as exhibited by Figure 5, the weaker expressing a particular gene is, the higher the tendency of the  $\log_{10}(\text{expression ratio})$  of the gene from two nominal repeats of an experiment to deviate from zero. Thus, the low-abundance (weakly expressing and hence low-intensity hybridization) genes will tend to occupy the tails of the distribution of expression ratios (i.e. deviate from zero in accordance with Figure 5) more often than the higher-abundance genes.

On page 19, line 11, of the specification, please replace the paragraph beginning “From Figure 5 it is evident” with the following amended paragraph:

From Figure 5 it is evident that the contour lines follow the error envelope. The value of  $d$  is proportional to the number of contours that a particular measurement falls away from  $\log(\text{Expression Ratio}) = 0$ . Thus the errors are distributed with respect to the contours similarly at low and at high intensity, and  $d$  has the desired property. One advantage of plotting contour lines is that the amount of error associated with each cellular constituent measured on the microarray can be calculated based on information derived from the variance of all the cellular constituents on the microarray across a plurality of measurements. Thus, by using grid lines as plotted in Figure 5, the significance of any deviation between  $X_i$  and  $Y_i$ , in a two-color fluorescent probe hybridization experiment, where  $i$  is a particular cellular constituent, will be placed in the context of the entire error envelope using an equation such as the denominator of equation 7. This provides an intensity independent method for determining the reliability of ~~measurement~~ measurements made of particular cellular constituents in microarray experiments including two-fluorophore or single-fluorophore experiments.

On page 21, line 1, of the specification, please replace the paragraph beginning “In practice the individual errors” with the following amended paragraph:

In practice the individual errors,  $\sigma_i^2$ , are themselves uncertain. Inspection of control experiments such as Figure 5 indicates the rough distribution of errors, but do not indicate whether individual genes at a particular intensity tend to have larger errors due to peculiarities of their RNA extraction or even biological function in the cell. Thus a better estimate of the error in the weighted mean is obtained by adding a component to equation (9) that accounts for scatter in the repeated measurements. If we denote the observed standard deviation for gene  $j$  as  $s_j$ , the error in the mean may be described as:

$$\sigma_x^2 = \frac{1}{N} \left[ \left( \sqrt{\frac{1}{\sum_i \frac{1}{\sigma_i^2}}} \right) + (N-1) * s_j \right] \quad (10)$$

where N is the number of repeated measurements. Equation (10) transitions from Equation (9) to the value of the observed scatter,  $s_j$ , as the number of repeats, N, becomes large. Note that  $s_j$  is calculated according to traditional statistical methods, such that

$$s_j \equiv \frac{1}{N-1} \sum_i (x_i - \bar{x})^2 \quad (11)$$

where N is the number of measurements,  $x_i$  are individual measurements of the intensity of gene j in a particular microarray experiment and  $\bar{x}$  is the sample mean of the individual measurements. See e.g. equation 2-10 in “Data Reduction and Error Analysis for the Physical Sciences”, Sciences,<sup>2</sup> *supra*, where  $s_j = \sigma^2$ . An estimate of the error of the mean,  $\bar{x}$ , as described by equation (10) is necessary because, equations such as (11) require a large number of nominal repeats (N) in order to be a true reflection of error. Estimates of error based on equation (9) do not take into consideration the errors that particular ~~measurement~~ measurements are susceptible to, as illustrated in Figure 1, ~~and~~ as well as gene specific anomalies. One skilled in the art will note that other equations that accomplish the transition from equation (9) to equation (10) are possible.

On page 37, line 31, of the specification, please replace the paragraph beginning “PCR products containing” with the following amended paragraph:

PCR products containing common 5' and 3' sequences (Research Genetics) were used as templates with amino-modified forward primer and unmodified reverse primers to PCR amplify 6065 ORFs from the *S. ~~cerevisiae~~ cerevisiae* genome. First pass success rate was 94%. Amplification reactions that gave products of unexpected sizes were excluded from subsequent analysis. ORFs that could not be amplified from purchased templates were amplified from genomic DNA. DNA samples from 100 µl reactions were isopropanol precipitated, resuspended in water, brought to 3x SSC in a total volume of 15 µl, and transferred to 384-well microtiter plates (Genetix). PCR products were spotted onto 1x3 inch polylysine-treated glass slides by a robot built according to specifications provided in Schena *et al.*, *supra*; DeRisi *et al.*, 1996, Discovery and analysis of inflammatory disease-related genes using micorarrays, PNAS USA, 94:2150-2155; and ~~DeResi~~ DeRisi *et al.*, (1997). After printing, slides were processed following published protocols. See ~~DeResi~~ DeRisi *et al.*, (1997).